

TJHSST Computer Systems Lab Senior
Research Project
Statistical Analysis of Mouse Gut Microbiota
2009-2010

Alex Tran

April 6, 2010

Abstract

The mouse gut microbial community is the population of bacteria, which inhabit the digestive track of a given population of mice. The ability to understand the microbial community has implications, which extend beyond mice and into factors of human obesity and the dietary needs of the human body. Computational genomics is an emerging field, which blends both biology and computer science, to analyze genes. This study seeks to utilize several methods emerging within computational genomics to analyze the gut microbiome of mice to observe the effects of varying diets on the gut microbial community, as well as create a universal and user friendly tool for researchers to utilize when studying any gut microbiome. The applications of this can extend past studying the gut microbiome, but can also be applied to any taxonomic group being counted for analysis, however several tools found in the software are exclusively beneficial to the analysis of gut microbiota specifically those found in mice.

Keywords:Computational Genomics, Genomics, Gut Microbiota, Statistics, Microbiota

1 Introduction

1.1 Scope of Study

To use genomics to observe and document the gut microbial community in the fecal matter of mice during all stages of analysis. From the sequencing of genetic material to the processing of the microbial community. To create an all encompassing suite to streamline the analysis of any microbial and community.

1.2 Expected results

To be able to present data from the microbial community in a way that is user friendly, efficient, and able to represent the data in an accurate way. Implemented Gui utilization that will display statistical information on each run, which can be used as visuals in a report, and to be able to do a full run from the ground up given 16sRNA sequencing data that can accurately and efficiently identify individual microbiota within the community. The end result of which will be an all comprehensive application that will be able to take generated data, and give the user a full evaluation of these data.

2 Literature Review

Other research has observed that general imbalances inside the gut microbiome can lead to differences in tendencies towards Obesity in mice(classified as OB/OB)(Turnbaugh, Gordon).It has also observed that within various populations of mice the microbiome appears to be so diverse that there is little known correlation between an individual and the population of their microbiome, however it was also shown that though the population in the microbiome is diverse there are tendencies based on the diet of an individual mouse for an increased representation of certain strains of bacteria to be present within the microbial community (Turnbaugh, et. al.)

Though General imbalances have been observed in the gut microbiome of obese micethere is still not enough information or analysis to say pinpoint what specific species of bacteria can lead to the tendency for mice to be obese.(Turnbaugh, Gordon) This lack of information is largely due to the

inefficiency and inaccurate nature of current computational 16s sequencing techniques (Turnbaugh et. al.)

To combat these inefficiencies a recent study attempted to use computational genomics to extrapolate microRNA from a "model organism" such as the *Drosophila* genus of flies and predict the structure of certain genes in other organisms of similar structure and species. To do this an algorithm was developed that would take microRNA archive it, and use it to predict the structure of certain marker genes that should be present in the following iteration. The hope is that if such an algorithm that can reliably predict the structure of genes given certain restraints could be created it would save a large sum of money on expensive DNA testing, and vastly improve the speed at which research can be done. (He Tao, Li Fei, Gu Jin, Li Ruiqiang, Li Fei)

3 Development

The program's end result consisted of several parts, which could be created individually in any order. Due to possible problems with time constraints it was been written in reverse depth from the most superficial part to the most abstract. This will allow for user friendly verification as well as a guaranteed presentable result if time had become a problem. It was primarily programmed in Python. Coding began by generating a simple taxonomic tree (originally defaulting to the RDP taxonomic tree,) and filling it with data points. It then moved onto generating a histogram of these data using Tkinter as a base. After this simple histogram was generated specific file loading, and a variety of taxonomic trees were made available to the user to allow the ability of the application to extend past only gut microbiota in mice to any accepted use of taxonomic OTU counting. Programming then moved on to allow for mouse movement to allow for more user friendliness, and to aid between two varying samples.

Following these simple processes had been implemented a system, which maintained the freedom to use any taxonomic tree desired, but also allowed the program to default to the most commonly used tree, and for the user to reopen recently used files by utilizing the cPickle archiving system for python. To do this it became necessary to create a hashmap, which would keep track of the number of times each tree was used, which the program

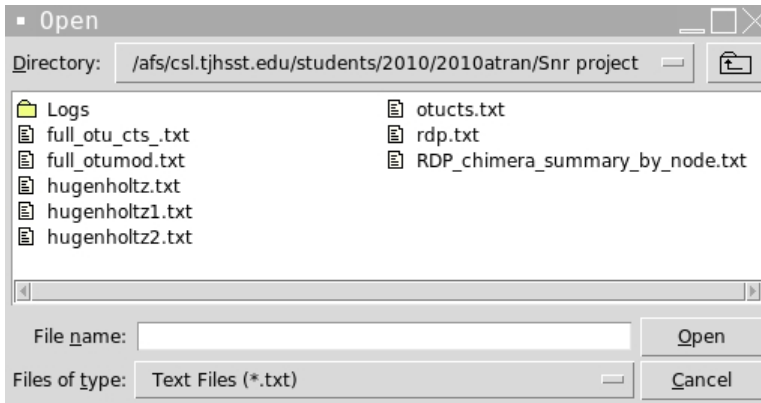


Figure 1: Screenshot of the file browser set to find .txt files

would default to if the user did not specify a taxonomic tree, as well as allow the user to call up the 10 most recently used files. The cPickle archiving system is not dissimilar to the way in which Microsoft Office's archiving system works in the way it uses a shared cache file between various instances of the program that is read when the program is opened, and edited upon closing the program. To test and verify results several premade computational genomics tools that each perform a part of what the finalized project does as one universal tool were used.

4 Quality assessment

Input data and 16sRNA sequences are to be received from the Gordon Lab at Washington University in St. Louis. The main goal of this project was to provide a user friendly OTU analysis suite, which could be used and interpreted by anyone familiar with genetics, but not necessarily with computer programming or use. To achieve this goal a large portion of development time and assessment was done in the minor debugging stage to present the user with the friendliest interface possible. To input a file the user is presented with the self explanatory file browser pictured in Figure.1 When a user inputs a file they are presented with the screen seen in Figure.2.

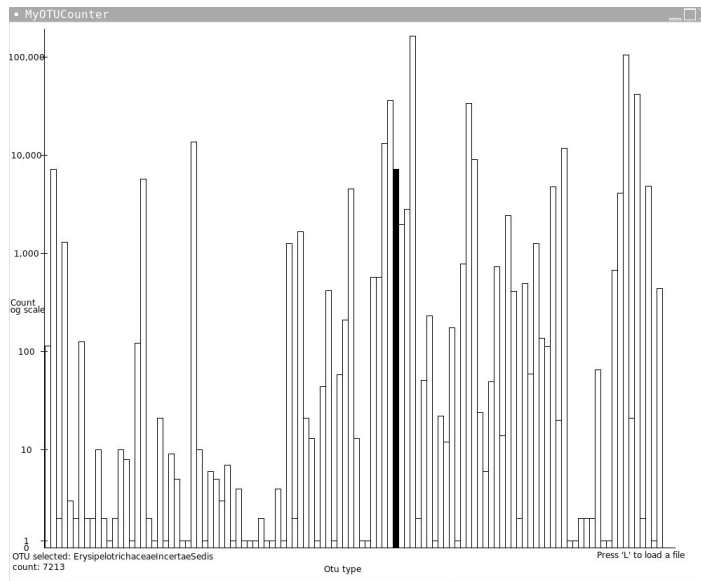


Figure 2: Screenshot of a loaded file, with a data point selected

5 Expected Results

To be able to present data from the microbial community in a way that is user friendly, efficient, and able to represent the data in an accurate way. To achieve this a gui was utilized that allows the display of statistical information on each run in a way that allows the user to more easily gather information from a sample than has ever before been possible. As a result of this analysis of data gathered from the gut microbiome can be done vastly faster than other archaic forms of analysis previously used. It cuts down on the amount of fact checking, and manual calculation that was previously necessary when doing analysis of this nature.

The results from this study are potentially very helpful to the field of genomics, because tools for analysis of gut microbial communities are a field that is in its early stages of development, and at this point most tools which are planned to be implemented into the end result application are only available separately, and very expensive.

It is also of note that while many of these tools are available in expen-

sive suites most of these require an in depth knowledge of both the field of genetics as well as a firm grasp on scripting in PERL or another comparable programming language, or are programmed to handle a very specific set of data. The goal of this project was to step past these former suites and provide a simple universal tool which could be used across the field or even apply to other fields related only in the need for adherence to a taxonomic tree. The goal was for the program to apply to any situation where counting of a taxonomic group under a given set of criteria is possible, which could vary from the previously stated examining of gut microbiota in mouse fecal matter to the number of a specific species of lemur in Madagascar. The program was therefore refined for gut microbiota, but given the freedom to adapt to many situations.

6 Discussion

The program in its early stages has proven to be useful in the Gordon Lab for which it was developed, and has also found use outside of the lab using taxonomic trees which were not tested during its development, which is a good sign for the universality of the program, however much could be done to improve upon it in its current state. The program has only served as a functional success largely due to the lack of a professional and user friendly suite similar to the one that was designed, which would likely contain a more complete set of analysis tools. Many features which a similar future program might offer would include the ability to compile and compare data from a number of test subjects, which in the current iteration must be done manually, and perhaps include a more from the ground up approach, which could not be completed due to a lack of access to the raw data by a claim from the company from which the Gordon Lab receives its analysis. Computational genomics is still in its infancy, and the tools which this programs provides, while not complete serve more as a proof of concept that temporarily fill the void where a lack of user friendly software currently exists.

7 Bibliography

Peter J. Turnbaugh and Jeffrey I. Gordon, "The core gut microbiome, energy balance, and obesity." Center for Genome Sciences, Washington University

School of Medicine, St. Louis MO, 2009.

Peter J. Turnbaugh, Ruth E. Ley, Micah Hamady, Claire M. Fraser-Liggett, Rob Knight Jeffrey I. Gordon, "The Human Microbiome Project." Center for Genome Sciences, Washington University School of Medicine, St. Louis MO, 2009.

He Tao, Li Fei, Gu Jin, Li Ruiqiang, Li Fei, "Computational Identification of 99 Insect MicroRNAs Using Comparative Genomics." Beijing Genomics Institute, Chienese Academy of Sciences, Beijing China, 2008.